



7/25/2023

Regression Analysis

GROUP-C

Tejaswi Mahapatra(21BECE46)
SILICON INSTITUTE OF TECHNOLOGY

Data Set Information: LC50 data, which is the concentration that causes death in 50% of test fish over a test duration of 96 hours, was used as a model response.

The data comprised 6 molecular descriptors:

MLOGP (molecular properties),
CIC0 (information indices),
GATS1i (2D autocorrelations),
NdssC (atom-type counts),
NdsCH ((atom-type counts),
SM1_Dz(Z) (2D matrix-based descriptors).

Details can be found in the quoted reference: M. Cassotti, D. Ballabio, R. Todeschini, V. Consonni.

Attribute Information: 6 molecular descriptors and 1 quantitative experimental response:

- 1) CIC0
- 2) SM1_Dz(Z) : 0 means missing value
- 3) GATS1i
- 4) NdsCH
- 5) NdssC
- 6) MLOGP
- 7) quantitative response, LC50 [-LOG(mol/L)]

Objective:

Prepare the dataset for analysis, using methods learnt and researched independently, including but not limited to Missing Value Imputation, Outlier Detection etc.

Find out your observations on different variables using descriptive statistics, Visualization etc. Report if there is any pattern present in the data.

Take LC50 [-LOG(mol/L)] as the target variable and fit different models for regression using other variables present in the data. Optimize the model parameters and find the best performing model. Compare all models used, using various performance metrics. Provide inferences to your findings.

Approach towards solving the business problem in QSAR_FISH_TOXICITY dataset:

- **Data Preprocessing & Exploration**
- **Feature Engineering & Transformation**
- **Feature Selection**
- **Data Scaling**
- **Model selection & hyperparameter tuning**
- **Model training & Evaluation**
- **Compare the result**
- **Finalize the model**
- **Interpretations on the model**
- **Ensemble methods(& interpretation for future predictions)**

Model Analysis:

Missing value treatment of target variable:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats.mstats import winsorize
from sklearn.model_selection import
train_test_split, KFold
from sklearn.linear_model import LinearRegression,
Ridge, Lasso, LogisticRegression
from sklearn.svm import SVR
from sklearn.tree import DecisionTreeRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor,
GradientBoostingRegressor
from sklearn.metrics import mean_squared_error,
mean_absolute_error, r2_score
from sklearn.model_selection import GridSearchCV
```

```

from fancyimpute import IterativeImputer
from sklearn.preprocessing import StandardScaler,
OneHotEncoder
from sklearn.decomposition import PCA
from scipy.stats import boxcox
from sklearn.metrics import make_scorer
from sklearn.linear_model import BayesianRidge
from xgboost import XGBRegressor
from lightgbm import LGBMRegressor
from sklearn.gaussian_process import
GaussianProcessRegressor
from sklearn.model_selection import cross_val_score

# Load the dataset
df = pd.read_csv("qsar_fish_toxicity.csv")

# Separate the target variable from the features
X = df.drop(columns=['LC50 [-LOG(mol/L)]'])
y = df['LC50 [-LOG(mol/L)]']
print(y.info())
# Impute missing values in the target variable using
IterativeImputer
y_imputer = IterativeImputer(verbose=False)
y_imputed =
y_imputer.fit_transform(y.values.reshape(-1, 1))

# Convert the imputed data back to a 1D Pandas Series
y_imputed = pd.Series(y_imputed.squeeze(), name='LC50
[-LOG(mol/L)]')
print(y_imputed.info())

```

- ✓ NECESSARY LIBRARIES WERE IMPORTED
- ✓ THE DATASET(.CSV) WAS LOADED INTO THE VARIABLE 'DF'
- ✓ THE TARGET VARIABLE WAS SEPARATED FROM THE FEATURES FOR THE UNDERLYING REASONS:

- SEPARATING X AND Y CLEARLY ESTABLISHES THE PREDICTION TASK.
- BY PROVIDING DISTINCT X AND Y , THE MODEL CAN UNDERSTAND THE MAPPING BETWEEN THE INPUT AND OUTPUT.
- THE MODEL'S PERFORMANCE IS ASSESSED BY MAKING PREDICTIONS ON NEW DATA USING X AND COMPARING THE PREDICTED Y WITH THE TRUE Y . THIS EVALUATION DETERMINES HOW WELL THE MODEL GENERALIZES TO UNSEEN DATA.
- ✓ UPON PRINTING `Y.INFO()` WE FOUND THAT THERE WERE MISSING VALUES IN THE TARGET COLUMN THAT NEEDED TO BE HANDLED (ON THE BASIS OF 908 ENTRIES OF THE DATASET)
- WE HANDLED THE MISSING VALUES IN THE TARGET COLUMN BY ITERATIVE IMPUTER: A MULTIPLE IMPUTATION TECHNIQUE.
- IT IS A MORE ACCURATE WAY TO IMPUTE MISSING VALUES THAN OTHER METHODS, SUCH AS MEAN IMPUTATION AND MEDIAN IMPUTATION. THIS IS BECAUSE MULTIPLE IMPUTATION TAKES INTO THE CORRELATION BETWEEN THE MISSING VALUES AND THE OBSERVED VALUES (SINCE WE HAVE LIMITED DOMAIN KNOWLEDGE, LOSING A LOT OF DATA OR COMPLETELY ERADICATING THE POSSIBILITY OF CORRELATION BETWEEN THE VALUES IS A SPECULATIVE CONJECTURE. HENCE, TO ME ITERATIVE IMPUTER LOOKED LIKE THE BEST POSSIBLE WAY TO HANDLE THE MISSING VALUES.)
- ✓ WE CONVERTED THE IMPUTED DATA BACK TO A 1D PANDAS SERIES IN ORDER TO BE ABLE TO PRINT THE INFO AND CHECK FOR THE HANDLED MISSING VALUES. HERE WE FOUND THAT THE MISSING VALUES WERE HANDLED SUCCESSFULLY. (908 ENTRIES: 908 NON-NULL)

File - Project

```
1 C:\Users\tejas\PycharmProjects\pythonProject\venv\Scripts\python.exe C:\Users\
  tejas\PycharmProjects\pythonProject\START\Project.py
2 <class 'pandas.core.series.Series'>
3 RangeIndex: 908 entries, 0 to 907
4 Series name: LC50 [-LOG(mol/L)]
5 Non-Null Count  Dtype
6 -----
7 906 non-null    float64
8 dtypes: float64(1)
9 memory usage: 7.2 KB
10 None
11 <class 'pandas.core.series.Series'>
12 RangeIndex: 908 entries, 0 to 907
13 Series name: LC50 [-LOG(mol/L)]
14 Non-Null Count  Dtype
15 -----
16 908 non-null    float64
17 dtypes: float64(1)
18 memory usage: 7.2 KB
19 None
```

Missing value treatment of feature variables & visualization :

```
# Print the information about the dataset
print(X.info())

# Plot a pair plot to see the relationships between
the different variables
sns.pairplot(X)
plt.show()

# Plot a heatmap to show the missing values in the
dataset
plt.figure(figsize=(8, 6))
sns.heatmap(X.isnull(), yticklabels=False,
cbar=False, cmap='viridis')
plt.title('Missing Value Heatmap')
plt.show()

# Calculate the skewness of each column
skewness = X.skew()
print(skewness)

# Check the distribution of the data in each column
for column in X.columns:
    print(column, X[column].describe())

# Calculate the missing-value percentage
missing_value_percent = X.isna().sum() / len(X) * 100
print("Missing Value Percentage:")
print(missing_value_percent)

# Impute the missing values using IterativeImputer
imputed_data =
IterativeImputer(verbose=False).fit_transform(X)

# Convert the imputed data to a Pandas DataFrame
imputed_df = pd.DataFrame(imputed_data,
columns=X.columns)
print(imputed_df.info())
```

- FIRST WE PRINTED THE INFORMATION ABOUT THE DATASET (THIS GIVES US THE NON-NULL VALUE COUNT)
- THEN WE PLOTTED THE PAIR-PLOT TO SEE THE RELATIONSHIP BETWEEN THE DIFFERENT VARIABLES (OPTED FOR A PAIR-PLOT AS IT SHOWS US ALL THE DEPENDENCIES AND RELATIONSHIPS IN ONE RUN)
 - **PAIR-PLOT OBSERVATION :**
 - WE CAN AVOID INCLUDING CICO AS IT IS JUST AN INFORMATION INDICE.
 - THE CLUSTERS APPEAR TO BE DISTRIBUTED RANDOMLY, AND THERE IS NO CLEAR PATTERN BETWEEN THE TWO FEATURES
 - THIS DOESN'T TELL MUCH ABOUT THE DATA, WE NEED TO PREPROCESS IT AND HANDLE THE DATA IN ORDER TO BE ABLE TO PREDICT, ASSESS AND CREATE THE BEST FIT FOR THE MODEL.
- PLOTTED A HEATMAP TO SHOW THE MISSING VALUES IN THE DATASET AS HEATMAPS PROVIDE AN INTUITIVE AND VISUAL WAY TO QUICKLY IDENTIFY PATTERNS AND AREAS WITH MISSING DATA. MISSING VALUES ARE OFTEN REPRESENTED BY A DISTINCT COLOUR, MAKING IT EASY TO SPOT REGIONS WHERE THE DATA IS INCOMPLETE.
 - **HEATMAP OBSERVATION:**
 - **THE COLUMNS WITH MOST MISSING VALUES:**
 - **SM1-DZ(z), CICO, GAZS1, NDSCH**
 - **NO MISSING VALUES IN NDSSC**
 - **FEWEST MISSING VALUE IN MLOGP**
- FURTHER ON, THE SKEWNESS WAS CALCULATED TO CHECK FOR THE POSITIVE OR NEGATIVE SKEWNESS AND CHOOSE A BETTER MISSING VALUE IMPUTATION METHOD.

- KNOWING THE SKEWNESS ALSO HELPED ME DECIDE IF ANY DATA TRANSFORMATION WAS NEEDED (AS SKEWNESS SHOWS US THE DISTRIBUTION)
- AT THIS POINT I ALREADY HAD THE INFORMATION THAT THE DATA IS SKEWED, SO I USED THE .DESCRIBE() METHOD TO CHECK THE DISTRIBUTION OF THE DATA IN EACH COLUMN AS CALCULATING MEASURES LIKE MEAN, MEDIAN AND STANDARD DEVIATION WOULD GIVE ME AN IDEA ABOUT THE DATA'S CENTRAL TENDENCY AND SPREAD, WHICH IS VERY INFORMATIVE TO CHOOSE BETTER DATA PREPROCESSING METHODS.
- FURTHER ON I CALCULATED THE MISSING VALUE PERCENTAGE, THIS IS AN ESSENTIAL STEP, AS IT IS HERE THAT WE GET TO KNOW HOW MUCH % OF OUR DATA IS MISSING AND WE CAN BETTER HANDLE THE MISSING VALUES FROM THIS INFORMATION.
 - FROM THIS WE GOT TO KNOW THAT ABOUT 22% OF THE DATA WAS MISSING. NOW IF WE DIDN'T IMPUTE, WE'D HAVE LOST A LOT OF DATA. ON TOP OF THAT THE MISSING VALUES WERE NOT MISSING AT RANDOM, THIS MEANS THE MISSING VALUES COULD BE CORRELATED WITH THE OBSERVED VALUES. IF WE DIDN'T IMPUTE, IT COULD'VE BIASED THE RESULT OF MY ANALYSIS.
 - CONSIDERING THE DISTRIBUTION (WHICH CLEARLY WASN'T UNIFORM), AND THE INABILITY TO CLEARLY ESTABLISH THE CORRELATION BETWEEN THE MISSING VALUES AND OBSERVED VALUES I FOUND ITERATIVE IMPUTER OR RATHER MULTIPLE IMPUTATION TO BE A TECHNIQUE BETTER SUITED TO PERFORM THE IMPUTATION AS IT TAKES INTO ACCOUNT THE CORRELATION BETWEEN THE MISSING AND OBSERVED VALUES AND GIVES A

MORE ACCURATE RESULT IN COMPARISON TO OTHER IMPUTATION METHODS.

- THEN THE IMPUTED DATA WAS CONVERTED TO A PANDAS DATA FRAME AND THE INFORMATION ON THE DATA WAS PRINTED WHICH CLEARLY SHOWED ALL THE MISSING VALUES WERE IMPUTED SUCCESSFULLY.(908 ENTRIES: 908 NON-NULL)

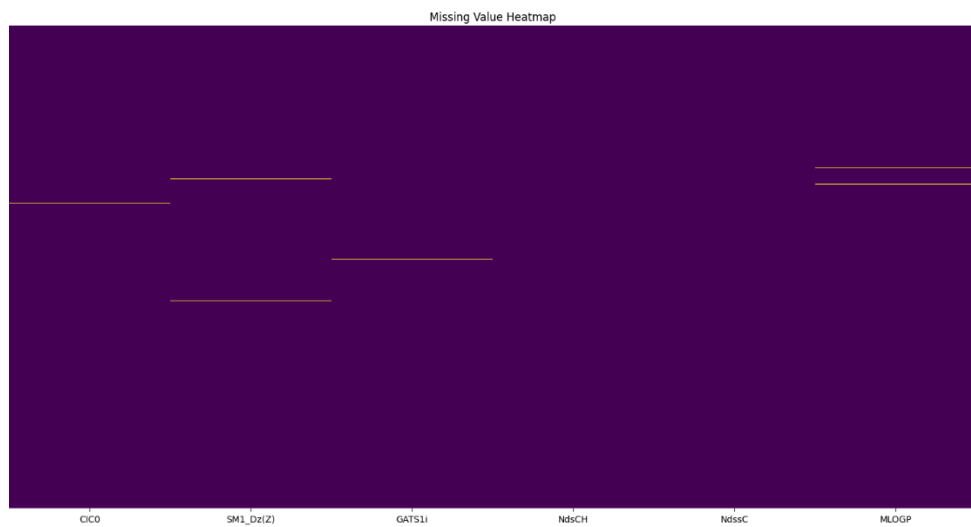
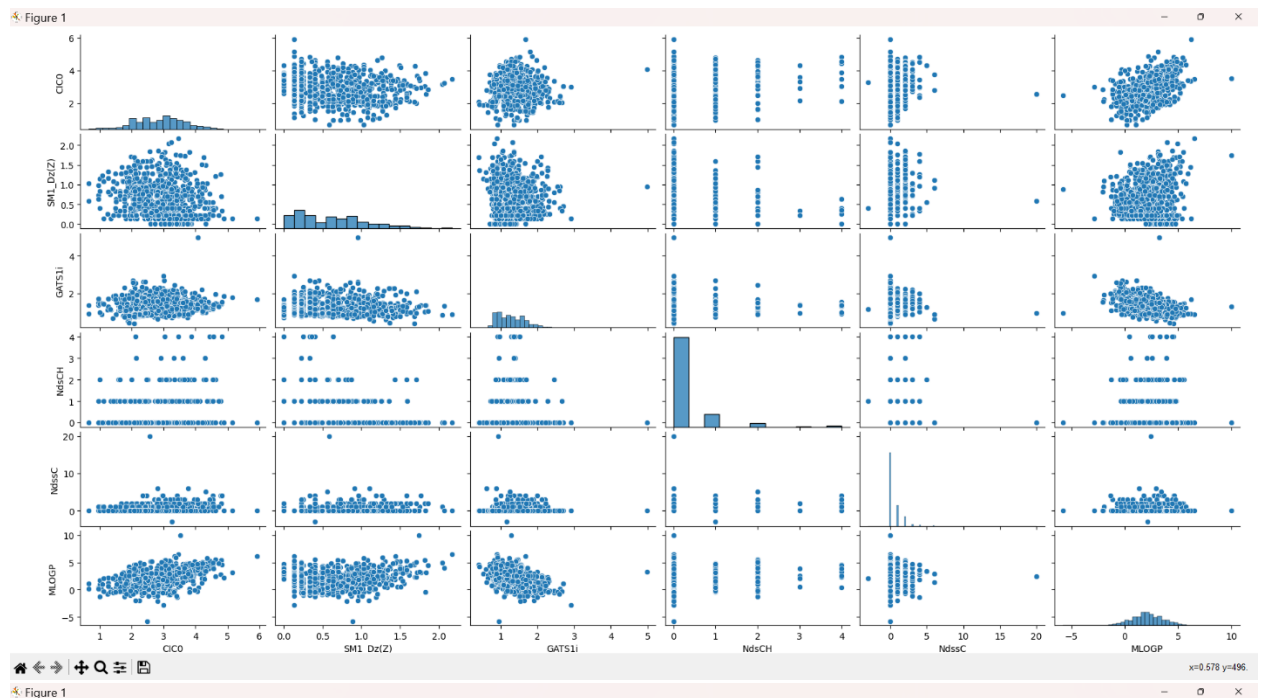
```

20 <class 'pandas.core.frame.DataFrame'>
21 RangeIndex: 908 entries, 0 to 907
22 Data columns (total 6 columns):
23 #   Column      Non-Null Count  Dtype
24 ---  ---
25 0    CIC0        906 non-null    float64
26 1    SM1_Dz(Z)   906 non-null    float64
27 2    GATS1i      906 non-null    float64
28 3    NdsCH       907 non-null    float64
29 4    NdssC       908 non-null    int64
30 5    MLOGP       905 non-null    float64
31 dtypes: float64(5), int64(1)
32 memory usage: 42.7 KB
33 None
34 CIC0        0.045111
35 SM1_Dz(Z)   0.694484
36 GATS1i      1.391213
37 NdsCH       3.398560
38 NdssC       7.489497
39 MLOGP       -0.038095
40 dtype: float64
41 CIC0 count   906.000000
42 mean        2.898620
43 std         0.756221
44 min         0.667000
45 25%         2.348750
46 50%         2.934000
47 75%         3.407000
48 max         5.926000
49 Name: CIC0, dtype: float64
50 SM1_Dz(Z) count  906.000000
51 mean         0.628595
52 std         0.428775
53 min         0.000000
54 25%         0.223000
55 50%         0.570000
56 75%         0.896250
57 max         2.171000
58 Name: SM1_Dz(Z), dtype: float64

```

```
59 GATS11 count      906.0000000
60 mean              1.297135
61 std               0.412765
62 min              0.396000
63 25%              0.950250
64 50%              1.240500
65 75%              1.562750
66 max              4.980000
67 Name: GATS1i, dtype: float64
68 NdsCH count       907.0000000
69 mean              0.229327
70 std               0.605621
71 min              0.000000
72 25%              0.000000
73 50%              0.000000
74 75%              0.000000
75 max              4.000000
76 Name: NdsCH, dtype: float64
77 NdssC count       908.0000000
78 mean              0.504405
79 std               1.083596
80 min             -3.000000
81 25%              0.000000
82 50%              0.000000
83 75%              1.000000
84 max             20.000000
85 Name: NdssC, dtype: float64
86 MLOGP count       905.0000000
87 mean              2.113989
88 std               1.480015
89 min             -5.780000
90 25%              1.209000
91 50%              2.127000
92 75%              3.109000
93 max             10.000000
94 Name: MLOGP, dtype: float64
```

```
95 Missing Value Percentage:
96 CICO          0.220264
97 SM1_Dz(Z)     0.220264
98 GATS1i        0.220264
99 NdsCH          0.110132
100 NdssC         0.000000
101 MLOGP         0.330396
102 dtype: float64
103 <class 'pandas.core.frame.DataFrame'>
104 RangeIndex: 908 entries, 0 to 907
105 Data columns (total 6 columns):
106 #   Column      Non-Null Count  Dtype
107 ---  ---
108 0    CICO        908 non-null   float64
109 1    SM1_Dz(Z)   908 non-null   float64
110 2    GATS1i      908 non-null   float64
111 3    NdsCH       908 non-null   float64
112 4    NdssC       908 non-null   float64
113 5    MLOGP       908 non-null   float64
114 dtypes: float64(6)
115 memory usage: 42.7 KB
116 None
```



Handling outliers:

```
# boxplot before imputing outliers
plt.figure(figsize=(10, 6))
imputed_df.boxplot()
```

```
plt.title('Box Plot (Before Outlier Handling)')
plt.show()

# Perform winsorization on each column to handle outliers
winsorized_df = imputed_df.apply(lambda x:
winsorize(x, limits=[0.1, 0.1]))

# boxplot after imputing outliers
plt.figure(figsize=(10, 6))
winsorized_df.boxplot()
plt.title('Box Plot (After Outlier Handling)')
plt.show()
```

➤ I USED A BOXPLOT FOR IDENTIFICATION OF OUTLIERS AS BOXPLOTS DISPLAY THE DATA'S FIVE-NUMBER SUMMARY, WHICH INCLUDES THE MINIMUM, FIRST QUARTILE(Q1),MEDIAN, THIRD QUARTILE(Q3), AND MAXIMUM. ANY DATA POINTS FALLING BELOW THE LOWER WHISKER($Q1 - 1.5 * IQR$) OR ABOVE THE UPPER WHISKER($Q3 + 1.5 * IQR$) ARE CONSIDERED POTENTIAL OUTLIERS AND ARE SHOWN AS INDIVIDUAL POINTS OUTSIDE THE BOXPLOT. ALSO, OUTLIERS THAT LIE FAR AWAY FROM THE MEDIAN AND QUARTILES ARE STILL CLEARLY VISIBLE, MAKING THEM EASIER TO IDENTIFY. ALONG WITH THIS THERE ARE SEVERAL OTHER BENEFITS OF VISUALIZING USING BOXPLOT.

- SO FROM THE BOX PLOT WE NOTICED ON WHICH CATEGORIES THE OUTLIERS WERE MAXIMUM AND COULD COMPARE THE DIFFERENT CATEGORIES ON THE BASIS OF OUTLIERS PRESENT FROM THE DATA AND FIGURE OUT HOW TO TREAT THEM.
- I DIDN'T USE Z-SCORE OR IQR OR OTHER METHODS TO DROP THE OUTLIERS AS WE COULD ESSENTIALLY LOSE A LOT ON THE

POTENTIAL DATA IF WE DROP THEM AND ALSO THERE ARE CONS LIKE Z-SCORES AREN'T ALWAYS EFFECTIVE AT IDENTIFYING OUTLIERS IN SKEWED DISTRIBUTIONS AS SKEWED DISTRIBUTIONS HAVE A LONG-TAIL WHICH MEANS THERE ARE A LOT OF VALUES FAR FROM THE MEAN, ETC

- SO IN PLACE OF THAT, I HANDLED THE OUTLIERS BY A METHOD CALLED AS 'WINSORIZATION'. WINSORIZATION IS A TECHNIQUE WHERE EXTREME VALUES ARE REPLACED WITH VALUES AT A SPECIFIED PERCENTILE (E.G. 95TH OR 99TH PERCENTILE). THIS APPROACH REDUCES THE INFLUENCE OF OUTLIERS WITHOUT REMOVING DATA PTS. ENTIRELY.
 - IF I SET THE PERCENTILE VALUE TO 5 MEANING THE TOP 5% AND THE BOTTOM 5% OF THE DATA WILL BE REPLACED WITH THE VALUE AT THE 5TH PERCENTILE AND 95TH PERCENTILE RESPECTIVELY.
 - A LOWER PERCENTILE VALUE (E.G. 1,5,10) WILL BE MORE AGGRESSIVE IN CAPPING EXTREME VALUES.
 - IN THIS PARTICULAR CODE I'VE SPECIFIED IT TO THE 10TH PERCENTILE FROM TOP OR (90TH FROM BOTTOM).
- SO NOW THE DATA HAS BEEN IMPUTED BY WINSORIZATION AND HANDLED FOR FURTHER PREPROCESSING AND FITTING INTO MODELS.
- AFTER OUTLIER HANDLING, THE BOXPLOT HAS CONSIDERABLE CHANGES AS THE NUMBER OF OUTLIERS NOW SEEM TO VERY LESS (CLOSE TO NONE) AND HAS BEEN HANDLED PROPERLY.

Figure 1

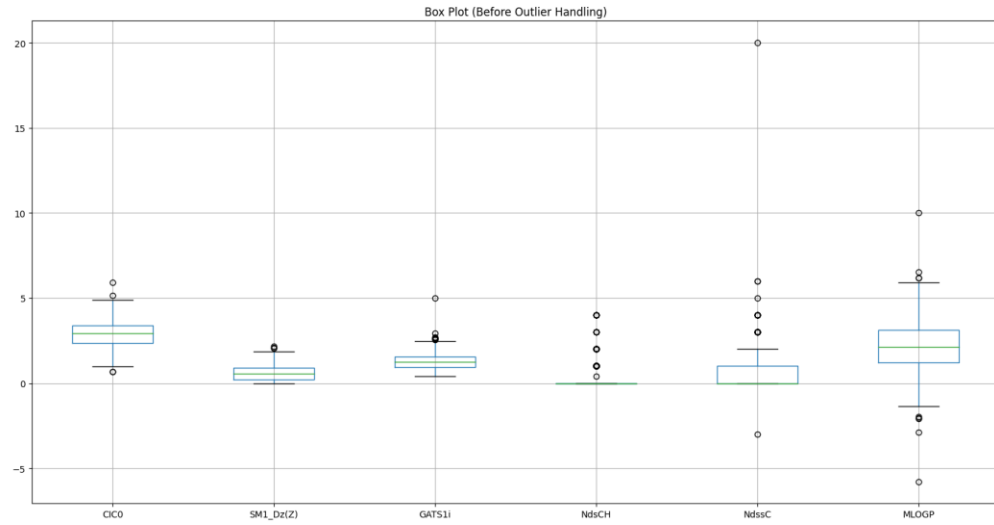


Figure 1

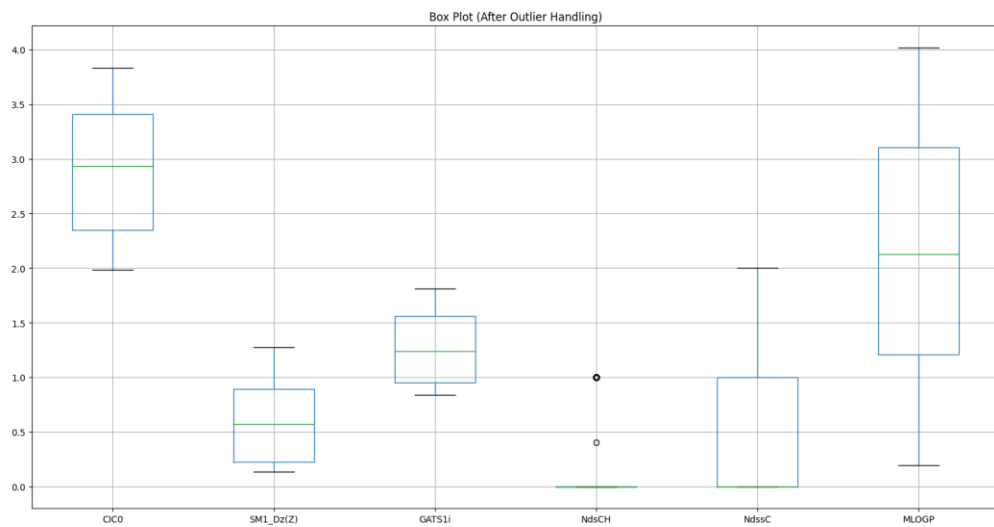


Figure 1

Data transformation:

```
# Calculate the correlation matrix before box-cox
correlation_matrix_before_boxcox =
winsorized_df.corr()
```

```
# Visualize the correlation matrix before box-cox
using a heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix_before_boxcox,
annot=True, cmap='coolwarm', fmt=".2f", center=0)
plt.title('Correlation Matrix (Before Box-Cox
Transformation)')
plt.show()

# skewness before box-cox
skewness_before_boxcox = winsorized_df.skew()
print("Skewness Before Box-Cox Transformation:")
print(skewness_before_boxcox)

# Apply Box-Cox transformation on each feature
for column in winsorized_df.columns:
    winsorized_df[column], _ = boxcox(
        winsorized_df[column] + 0.001) # Add a small
constant to handle non-positive values

# Calculate the correlation matrix after box-cox
correlation_matrix_after_boxcox =
winsorized_df.corr()

# Visualize the correlation matrix after box-cox
using a heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix_after_boxcox,
annot=True, cmap='coolwarm', fmt=".2f", center=0)
plt.title('Correlation Matrix (After Box-Cox
Transformation)')
plt.show()

# skewness after box-cox
skewness_after_boxcox = winsorized_df.skew()
print("Skewness After Box-Cox Transformation:")
print(skewness_after_boxcox)
```

- BEFORE APPLYING BOX-COX TRANSFORMATION, IT IS ESSENTIAL TO EXAMINE THE DATA ,UNDERSTAND ITS CHARACTERISTICS, AND MAKE INFORMED DECISIONS BASED ON THE CORRELATION STRUCTURE AND SKEWNESS.
- THE BOX-COX TRANSFORMATION IS A POWERFUL TRANSFORMATION THAT OPTIMISES THE TRANSFORMATION PARAM(LAMBDA) TO MAKE THE DATA AS CLOSE TO NORMAL AS POSSIBLE.
- IN THE ABOVE CODE, CORRELATION MATRIX HAS BEEN VISUALIZED USING HEATMAP BEFORE BOX-COX AND AFTER BOX-COX, AS CORRELATION MATRIX HELPS US UNDERSTAND THE RELATIONSHIP BETWEEN DIFFERENT VARIABLES IN THE DATASET. IT SHOWS THE PAIRWISE CORRELATIONS BETWEEN ALL PAIRS OF VARIABLES, INDICATING HOW THEY'RE RELATED TO EACH OTHER(+VELY OR -VELY)
 - HEATMAP(BEFORE BOX-COX) OBSERVATION:
 - THERE ARE NO MAJOR STRONG POSITIVE OR NEGATIVE CORRELATIONS
 - STRONG POSITIVE CORRELATION IS SEEN BETWEEN MLOGP & CICO THAT IS 0.44 & STRONG NEGATIVE CORRELATION IS SEEN BETWEEN MLOGP & GATS1 THAT IS -0.42
 - SIMILARLY THERE ARE OTHER INSIGHTS THAT WE CAN OBSERVE ABOUT EACH CATEGORY'S CORRELATION AND COMPARE
- SKEWNESS MEASURES THE ASYMMETRY OF THE DATA DISTRIBUTION. POSITIVE SKEWNESS INDICATES A LONGER TAIL ON THE RIGHT, WHILE NEGATIVE SKEWNESS INDICATES A LONGER TAIL ON THE LEFT.
 - BASED UPON THE SKEWNESS METRICS , I DECIDED TO OPT FOR BOX-COX TRANSFORMATION(SINCE IT'S ONLY APPLICABLE FOR POSITIVE VALUES WE ADD A SMALL CONSTANT TO HANDLE THE

NON POSITIVE VALUES) ON ALL THE COLUMNS, FOR A MORE NORMAL DISTRIBUTION OF THE DATA AND THE CHANGE IN SKEWNESS BEFORE AND AFTER BOX-COX HAS BEEN TABULATED BELOW:

<u>Skewness Before Box-Cox Transformation:</u>		<u>Skewness After Box-Cox Transformation:</u>	
CIC0	-0.000856	CIC0	-0.068212
SM1_Dz(Z)	0.374971	SM1_Dz(Z)	-0.108152
GATS1i	0.221689	GATS1i	-0.010239
NdsCH	1.823879	NdsCH	1.816913
NdssC	1.289047	NdssC	0.793025
MLOGP	-0.017874	MLOGP	-0.280295

- WE CAN SEE NOW THAT AFTER BOX-COX TRANSFORMATION THE DATA IS FAIRLY NORMALLY DISTRIBUTED. BOX-COX IMPROVED THE LINEAR RELATIONSHIP BETWEEN THE TRANSFORMED VARIABLES SERVING ITS PURPOSE.
- HEATMAP (AFTER BOX-COX) OBSERVATION:
 - STRONGEST POSITIVE CORRELATION IS SEEN BETWEEN MLOGP & CICO THAT IS 0.43 BEFORE BOX-COX IT WAS 0.44
 - STRONGEST NEGATIVE CORRELATION IS SEEN BETWEEN MLOGP & GATS1i THAT IS -0.42

117 Skewness Before Box-Cox Transformation:

Page 2 of 4

File - Project

118	CIC0	-0.000856
119	SM1_Dz(Z)	0.374971
120	GATS1i	0.221689
121	NdsCH	1.823879
122	NdssC	1.289047
123	MLOGP	-0.017874
124	dtype: float64	
125	Skewness After Box-Cox Transformation:	
126	CIC0	-0.068212
127	SM1_Dz(Z)	-0.108152
128	GATS1i	-0.010239
129	NdsCH	1.816913
130	NdssC	0.793025
131	MLOGP	-0.280295
132	dtype: float64	

Figure 1

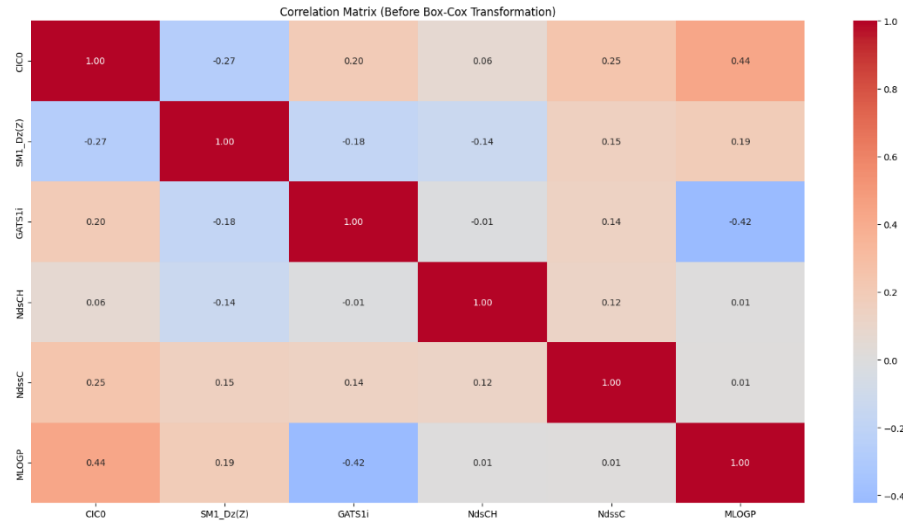


Figure 1

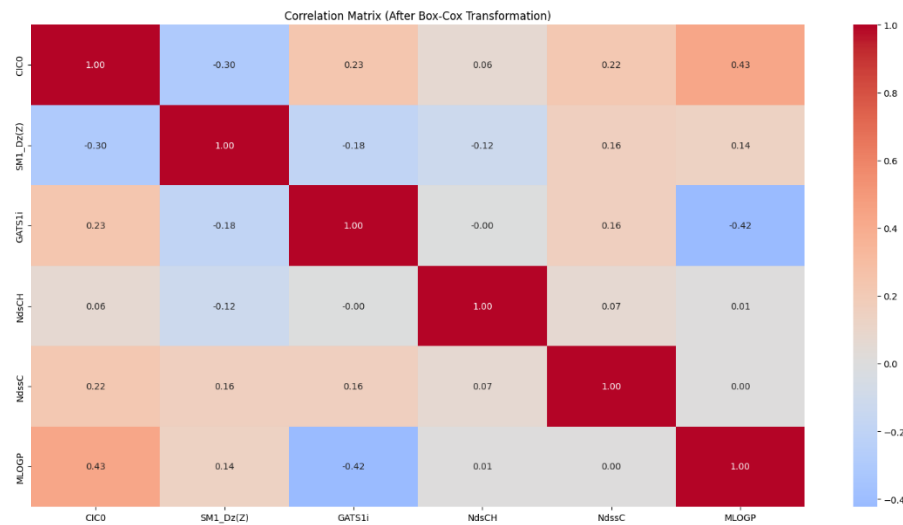


Figure 1

Scaling, Dimensionality reduction, encoding, descriptive statistics, visualization, and insights:

```
# Scale the data using StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(winsorized_df)
```

```

# Apply PCA for dimensionality reduction
# pca = PCA(n_components=5)
# X_pca = pca.fit_transform(X_scaled)

# One-hot encode the categorical features: "CIC0" and
"SM1_Dz(Z)"
encoder = OneHotEncoder(drop='first',
sparse_output=False)
X_encoded =
pd.DataFrame(encoder.fit_transform(winsorized_df[['CIC0', 'SM1_Dz(Z)']]))
X_encoded.columns =
encoder.get_feature_names_out(['CIC0', 'SM1_Dz(Z)'])

# Concatenate the one-hot encoded categorical columns
with the remaining features
X_final = pd.concat([X_encoded,
pd.DataFrame(X_scaled,
columns=winsorized_df.columns)], axis=1)

# Descriptive Statistics
print(X_final.describe())

# Create a scatter plot
plt.scatter(X_final['MLOGP'], y_imputed, label='Data
Points', color='blue', marker='o')
# Set plot labels and title
plt.xlabel('MLOGP')
plt.ylabel('LC50 [-LOG(mol/L)]')
plt.title('Scatter Plot Between Molecular prop(MLOGP)
& Target value(LC50)')
# Show the legend
plt.legend()
# Display the plot
plt.show()

# Pairwise Scatter Plots
features_to_plot = ['GATS1i', 'NdsCH', 'NdssC',
'MLOGP']

```

```
sns.pairplot(data=pd.concat([X_final[features_to_plot], y_imputed], axis=1), hue='LC50 [-LOG(mol/L)]')  
plt.show()
```

- FIRST, WE'VE SCALED THE DATA. SCALING IS USED FOR SEVERAL REASONS THAT INCLUDE SOLVING MAGNITUDE DIFFERENCES(MAKING A CONSTANT RANGE E.G. 0-1 OR 0-100,ETC), NORMALIZATION, CONVERGENCE AND EFFICIENCY ETC.
- SCALING IS ALSO BENEFICIAL FOR BETTER PERFORMANCE OF A LOT OF MODELS BUT IT'S SUBJECT TO THE DATASET THAT'S AVAILABLE & ITS ANALYSIS.
- IN THIS PARTICULAR CASE I'VE DONE SCALING TO MATCH THE MAGNITUDE AND SOLVE THE MAGNITUDE DIFFERENCE TO GET IT ON A CERTAIN RANGE SCALE.
- PCA (OR PRINCIPAL COMPONENT ANALYSIS)IS EFFECTIVE IF THERE ARE STRONG CORRELATIONS(+VE OR -VE).SINCE WE SAW THROUGH THE HEATMAP(BEFORE &AFTER BOX-COX) THERE ARE NO MAJOR STRONG CORRELATIONS THERE'S NO NEED TO USE DIMENSIONALITY REDUCTION TECHNIQUE LIKE PCA(IDENTIFIES A NEW SET OF UNCORRELATED VALUE THAT REPRESENTS THE MOST IMPORTANT PATTERNS OR DIRECTIONS OF VARIATION IN THE ORIGINAL DATA). HENCE, I'VE COMMENTED THIS (THE PCA)PART OUT.
- THE INFO DOESN'T SAY THERE ARE ANY CATEGORICAL VALUES IN THE DATASET BUT BASED ON MY UNDERSTANDING OF THE DATA AND THE REFERENCES AVAILABLE I CONSIDER CICO AND SM1_DZ(Z) TO BE CATEGORICAL VALUES. CICO(INFORMATION INDICES), SM1(2D MATRIX-BASED DESCRIPTORS) BOTH OF WHICH ARE METRICS WITH WHICH THE TOXIN CAN BE IDENTIFIED.

- **CICO**: CHEMICAL IDENTIFICATION CODE OF THE CHEMICAL. IT HAS FINITE NUMBER OF CATEGORIES.
- **SM1_DZ(Z)**: SHAPE OF THE MOLECULE, THAT IS FINITE NUMBER OF CATEGORIES.
- ONE-HOT ENCODING IS TYPICALLY USED FOR CATEGORICAL FEATURES TO CONVERT THEM INTO A NUMERICAL REPRESENTATION THAT CAN BE USED IN MACHINE LEARNING MODELS. FOR EXAMPLE, IF **CICO** CAN TAKE ON THE VALUES "A", "B", AND "C", THEN OHE WOULD CONVERT IT INTO THREE BINARY FEATURES: **CICO_A**, **CICO_B**, AND **CICO_C**. EACH OF THESE BINARY FEATURES WOULD BE 1 IF THE ORIGINAL **CICO** FEATURE WAS EQUAL TO THAT VALUE, AND 0 OTHERWISE.
- OHE IS DONE FOR CATEGORICAL FEATURES BECAUSE IT ALLOWS MACHINE LEARNING ALGORITHMS TO BETTER UNDERSTAND THE RELATIONSHIPS BETWEEN THESE FEATURES AND THE TARGET VARIABLE.
- FOR EXAMPLE, IF A MACHINE LEARNING ALGORITHM IS TRYING TO PREDICT THE TOXICITY OF A CHEMICAL, IT WOULD BE HELPFUL FOR THE ALGORITHM TO KNOW THAT THE CHEMICAL HAS A **CICO** VALUE OF "A". BY CONVERTING **CICO** INTO THREE BINARY FEATURES, THE ALGORITHM CAN LEARN THAT THE CHEMICAL IS MORE LIKELY TO BE TOXIC IF IT HAS A **CICO** VALUE OF "A".
- IN ADDITION, OHE CAN HELP TO PREVENT OVERFITTING. OVERFITTING OCCURS WHEN A MACHINE LEARNING ALGORITHM LEARNS THE TRAINING DATA TOO WELL AND IS NOT ABLE TO GENERALIZE TO NEW DATA. BY CONVERTING CATEGORICAL FEATURES INTO BINARY FEATURES, OHE CAN HELP TO REDUCE THE NUMBER OF PARAMETERS IN THE MACHINE LEARNING MODEL, WHICH CAN HELP TO PREVENT OVERFITTING.
- THEN WE CONCATENATED THE ONE-HOT ENCODED CATEGORICAL COLUMNS WITH THE REMAINING FEATURES.
- THEN WE OBSERVED THE STATISTICS OF THE NEW **X_FINAL** DATASET AFTER THE ENTIRE PREPROCESSING AND EXPLORATION.

- I DIDN'T DO FEATURE SELECTION AS WE DO NOT HAVE A LARGE NUMBER OF FEATURE, ELSE WE COULD'VE DONE LASSO OR UNIVARIATE FEATURE SELECTION ON THIS DATA. THERE'S ALSO MULTICOLLINEARITY HANDLING, HANDLING THE IMBALANCED DATA, DATA NORMALIZATION, HANDLING TIME-SERIES DATA ETC THAT COMES UNDER DATA CLEANING AND PREPROCESSING BUT FOR THE GIVEN DATASET THE ABOVE-MENTIONED TECHNIQUES ARE NOT IMPORTANT AND THE PRE-PROCESSING STEPS ALREADY DONE ARE SUFFICIENT.

- WE DID A SCATTER PLOT BETWEEN MOLECULAR PROP(MLOGP) & TARGET VALUE(LC50)

OBSERVATION:

THE SCATTER PLOT SUGGESTS THAT THEY ARE KIND OF LINEARLY RELATED BUT MLOGP CANNOT BE THE ONLY PARAMETER OR METRIC THAT CHANGES/ AFFECTS THE LC50. THE SCATTERPLOT SHOWS THAT THERE ARE SOME POINTS THAT FALL BELOW THE LINE OF BEST FIT. SO, NO EXACT OBSERVATION CAN BE MADE OUT OF THIS PARTICULAR SCATTER PLOT.

- SINCE INDIVIDUAL SCATTER PLOTS COULD TAKE A LOT OF TIME AND BE REALLY EXTENSIVE I THEN DID A PAIRWISE PLOT

OBSERVATION:

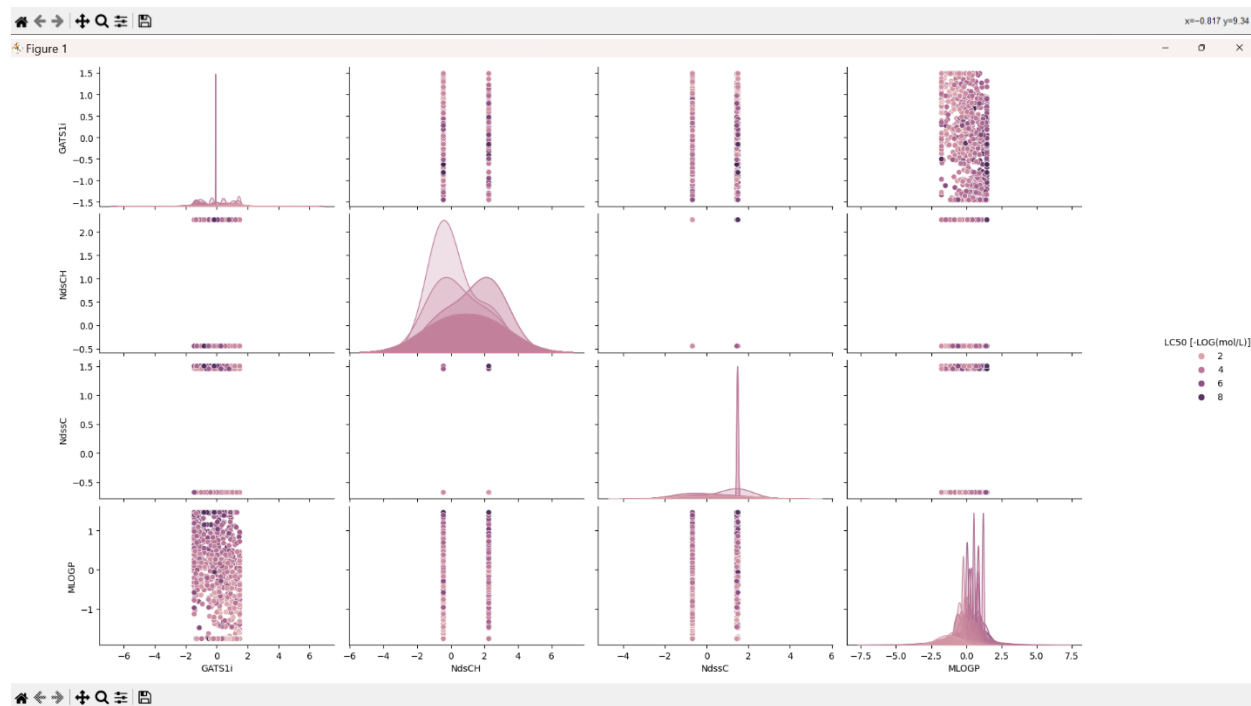
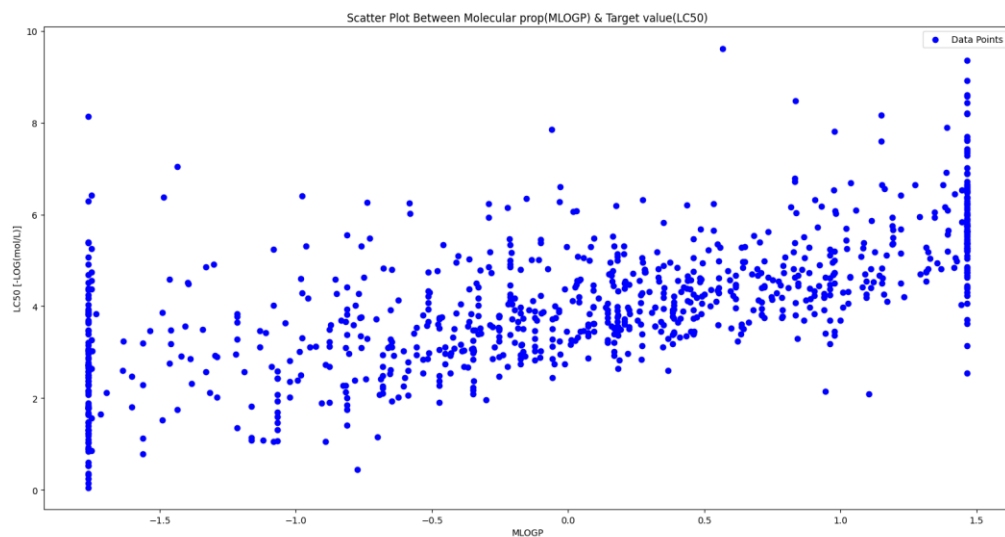
- GATS1i AND MLOGP ARE POSITIVELY CORRELATED, MEANING THAT AS THE VALUE OF GATS1i INCREASES, THE VALUE OF MLOGP ALSO TENDS TO INCREASE. THIS SUGGESTS THAT GATS1i AND MLOGP ARE BOTH MEASURES OF THE LIPOPHILICITY OF A MOLECULE, AND THAT THEY ARE BOTH CORRELATED WITH TOXICITY.
- NDSSC AND MLOGP ARE ALSO POSITIVELY CORRELATED, SUGGESTING THAT NDSSC IS ALSO A MEASURE OF THE LIPOPHILICITY OF A MOLECULE.
- NdsCH AND LC50 ARE NEGATIVELY CORRELATED, MEANING THAT AS THE VALUE OF NdsCH INCREASES, THE VALUE OF LC50 TENDS TO DECREASE. THIS SUGGESTS THAT NdsCH IS A MEASURE OF THE

NON-TOXICITY OF A MOLECULE, AND THAT IT IS INVERSELY
CORRELATED WITH TOXICITY.

✓ *A model could be trained to predict the toxicity of a molecule
based on its GATS1i, NdssC, and NdsCH values.*

133	CIC0_0.89515523512186	...	MLOGP
134	count	908.000000	9.080000e+02
135	mean	0.001101	5.869020e-17
136	std	0.033186	1.000551e+00
137	min	0.000000	-1.763288e+00
138	25%	0.000000	-6.799279e-01
139	50%	0.000000	9.232111e-02
140	75%	0.000000	8.297244e-01
141	max	1.000000	1.466201e+00
142			

Figure 1



Train Test split and choice of the regression model:

```
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test =
train_test_split(X_final, y_imputed, test_size=0.1,
random_state=42)
# Set up regression models and their hyperparameter
grids for tuning
```

```
models = {
    'Linear Regression': {
        'model': LinearRegression(),
        'params': {}
    },
    'Ridge Regression': {
        'model': Ridge(),
        'params': {
            'alpha': [0.1, 1.0, 10.0]
        }
    },
    'Lasso Regression': {
        'model': Lasso(),
        'params': {
            'alpha': [0.1, 1.0, 10.0]
        }
    },
    'SVR': {
        'model': SVR(),
        'params': {
            'C': [1, 10],
            'kernel': ['linear', 'rbf']
        }
    },
    'Random Forest Regressor': {
        'model': RandomForestRegressor(),
        'params': {
            'n_estimators': [50, 100, 200],
            'max_depth': [None, 5, 10]
        }
    },
    'Decision Tree Regressor': {
        'model': DecisionTreeRegressor(),
        'params': {
            'max_depth': [None, 5, 10]
        }
    },
    'XGBoost Regressor': {
        'model': XGBRegressor(),
```

```

        'params': {
            'learning_rate': [0.01, 0.1],
            'n_estimators': [100, 200],
            'max_depth': [3, 5]
        }
    },
    'KNN Regressor': {
        'model': KNeighborsRegressor(),
        'params': {
            'n_neighbors': [3, 5, 7]
        }
    },
    'Gaussian Process Regressor': {
        'model': GaussianProcessRegressor(),
        'params': {}
    },
    'Bayesian Ridge Regressor': {
        'model': BayesianRidge(),
        'params': {
            'alpha_1': [1e-6, 1e-5, 1e-4],
            'alpha_2': [1e-6, 1e-5, 1e-4],
            'lambda_1': [1e-6, 1e-5, 1e-4],
            'lambda_2': [1e-6, 1e-5, 1e-4]
        }
    },
},
}

```

- THE DATASET IS BEING SPLIT INTO TRAINING AND TESTING SETS. THE TRAINING SET WILL BE USED TO TRAIN THE REGRESSION MODELS, AND THE TESTING SET WILL BE USED TO EVALUATE THEIR PERFORMANCE.
- THE CODE SETS UP VARIOUS REGRESSION MODELS WITH DIFFERENT HYPERPARAMETERS, WHICH CAN BE USED FOR HYPERPARAMETER TUNING AND MODEL SELECTION BASED ON THEIR PERFORMANCE ON THE TRAINING AND TESTING DATASETS.
- BY CREATING 'MODELS' DICTIONARY I'VE ACHIEVED THIS IDEA.

- I ALSO OPTED FOR LIGHTGBM, QUANTILE, LOGISTIC, ELASTICNET, POLYNOMIAL, & PERCEPTRON ETC BUT ON TRYING THOSE I WASN'T GETTING AN IDEAL VALUE OR THE DESIRED METRICS. IN ORDER TO REDUCE THE NUMBER OF REGRESSORS USED I REMOVED THE ONES WITH ALMOST THE SAME OR WORSE PERFORMANCE METRIC AND PROCEEDED WITH THE FEW THAT I'VE MENTIONED ABOVE.

Optimisation of the model:

```
# Function to perform GridSearchCV and return the
best model and parameters
def find_best_model(model, params):
    grid_search = GridSearchCV(model, params,
scoring='neg_mean_squared_error', cv=5)
    grid_search.fit(X_final, y_imputed)
    best_model = grid_search.best_estimator_
    best_params = grid_search.best_params_
    return best_model, best_params

# Cross-validation for each model
for model_name, model_info in models.items():
    best_model, best_params =
find_best_model(model_info['model'],
model_info['params'])
    print(f"Best hyperparameters for
{model_name}: {best_params}")

    # Perform cross-validation on the entire
dataset
    cv_scores = cross_val_score(best_model,
X_final, y_imputed, cv=5,
scoring='neg_mean_squared_error')
    cv_rmse_scores = np.sqrt(-cv_scores)

    print(f"{model_name} - Cross-Validation RMSE:
{cv_rmse_scores.mean()}, Cross-Validation R2:
```

```
{cv_scores.mean()}")  
print("-----")
```

- I DEFINED A FUNCTION CALLED FIND_BEST_MODEL, WHICH PERFORMS HYPERPARAMETER TUNING USING GRIDSEARCHCV AND RETURNS THE BEST MODEL AND ITS CORRESPONDING BEST HYPERPARAMETERS. IT THEN APPLIES CROSS-VALIDATION ON EACH MODEL USING THE OPTIMIZED HYPERPARAMETERS AND EVALUATES THEIR PERFORMANCE USING ROOT MEAN SQUARED ERROR (RMSE) AND R-SQUARED (R2) METRICS.
- I WANTED TO BE SURE THAT I'D GET THE BEST POSSIBLE PERFORMANCE, AND GRIDSEARCHCV IS A MORE EXHAUSTIVE SEARCH ALGORITHM, WHICH MEANS THAT IT WILL TRY OUT MORE COMBINATIONS OF HYPERPARAMETERS. MAKING IT YIELD A BETTER TUNING.
- HELD OUT 10% OF THE TOTAL DATA. APPLIED 5-FOLD CV ON THE 90% DATA. ROBUST ENOUGH FOR CHECKING OVERFITTING.
- PRINTED THE SCORES.

Model Fitting, Training & Evaluation:

```
# Train the best model on the entire dataset  
best_model.fit(X_final, y_imputed)  
  
# Make predictions on the testing data  
y_pred = best_model.predict(X_final)  
  
# Evaluate the model's performance  
mse = mean_squared_error(y_imputed, y_pred)  
mae = mean_absolute_error(y_imputed, y_pred)  
r2 = r2_score(y_imputed, y_pred)  
adjusted_r2 = 1 - (1 - r2) * (len(y_imputed) - 1) /  
    (len(y_imputed) - X_final.shape[1] - 1)  
  
print(f"{model_name} - MSE: {mse}, MAE: {mae}, R-  
squared: {r2}, Adjusted R-squared: {adjusted_r2}")  
print("-----")
```


- AFTER OBTAINING THE BEST MODEL AND ITS OPTIMIZED HYPERPARAMETERS, THE BEST_MODEL (THE MODEL WITH THE BEST HYPERPARAMETERS FOUND THROUGH GRIDSEARCHCV) IS TRAINED ON THE ENTIRE DATASET X_FINAL AND Y_IMPUTED.
- THE TRAINED MODEL IS THEN USED TO MAKE PREDICTIONS ON THE SAME DATASET X_FINAL, WHICH WAS USED FOR TRAINING.
- SEVERAL PERFORMANCE METRICS ARE CALCULATED TO EVALUATE THE MODEL'S PERFORMANCE: MEAN SQUARED ERROR (MSE): IT MEASURES THE AVERAGE SQUARED DIFFERENCE BETWEEN PREDICTED VALUES AND ACTUAL TARGET VALUES. MEAN ABSOLUTE ERROR (MAE): IT MEASURES THE AVERAGE ABSOLUTE DIFFERENCE BETWEEN PREDICTED VALUES AND ACTUAL TARGET VALUES. R-SQUARED (R²): IT REPRESENTS THE PROPORTION OF VARIANCE IN THE TARGET VARIABLE EXPLAINED BY THE MODEL. IT IS A MEASURE OF HOW WELL THE MODEL FITS THE DATA. ADJUSTED R-SQUARED: IT IS A MODIFIED VERSION OF R-SQUARED THAT TAKES INTO ACCOUNT THE NUMBER OF FEATURES USED IN THE MODEL. IT PENALIZES THE INCLUSION OF IRRELEVANT FEATURES.
- THEN WE PRINTED THE CALCULATED MSE, MAE, R-SQUARED, AND ADJUSTED R-SQUARED FOR THE BEST MODEL.

COMPARISON & INTERPRETATION:

<u>REGRESSION MODEL</u>	<u>R-SQUARED SCORE</u>
<u>LINEAR REGRESSION</u>	<u>0.8313722680859739</u>
<u>RIDGE REGRESSION</u>	<u>0.6614223783722861</u>
<u>LASSO REGRESSION</u>	<u>0.5297899323478972</u>
<u>KNN REGRESSOR</u>	<u>0.7148848607010814</u>
<u>SVR</u>	<u>0.7830841904719806</u>
<u>RANDOM FOREST REGRESSOR</u>	<u>0.9344300181651313</u>
<u>GAUSSIAN PROCESS REGRESSOR</u>	<u>0.9841189000123955</u>

<u>BAYESIAN RIDGE REGRESSOR</u>	<u>0.6668329312859386</u>
<u>DECISION TREE REGRESSOR</u>	<u>0.6688594104076014</u>
<u>XGBOOST REGRESSOR</u>	<u>0.751346648231083</u>

- GAUSSIAN PROCESS REGRESSOR HAS THE BEST R-SQUARED SCORE (0.984) MEANING IT IS THE BEST POSSIBLE FIT FOR THE MODEL.
- LOWER CROSS-VALIDATION RMSE & MSE INDICATE BETTER PERFORMANCE IN TERMS OF PREDICTIONS ACCURACY.
- HIGHER CROSS VALIDATION R-SQUARED AND R-SQUARED VALUES CLOSE TO 1 INDICATE BETTER GOODNESS OF FIT.
- WE CAN DO THE FUTURE PREDICTIONS OF THE QSAR-FISH-TOXICITY TEST BASED ON THIS MODEL. FROM THE ENTIRE CODE ANALYSIS WE FOUND THAT THE FEATURES MOST INFLUENTIAL IN PREDICTING THE TARGET VALUE ARE GATS1I(MEASURE OF TOXICITY), NDSSC(MEASURE OF TOXICITY), AND NdsCH (MEASURE OF NON-TOXICITY)VALUES.
- ON THE BASIS OF EVALUATION AND PERFORMANCE METRIC WE NOW KNOW THAT THE MACHINE LEARNING MODEL MUST BE TRAINED ON GAUSSIAN PROCESS REGRESSOR FOR A BETTER FIT OF THE UNSEEN DATA.

Ensemble methods:

```
# Train the best models on the entire dataset
best_models = {
    'Random Forest Regressor':
RandomForestRegressor(n_estimators=100,
max_depth=None),
    'XGBoost Regressor':
XGBRegressor(learning_rate=0.1, n_estimators=200,
max_depth=5),
    'Gaussian Process Regressor':
GaussianProcessRegressor()
}
```

```
for model_name, model in best_models.items():
    model.fit(X_final, y_imputed)

# Make predictions using each individual model
predictions = {}
for model_name, model in best_models.items():
    predictions[model_name] = model.predict(X_final)

# Take the mean of the individual model predictions
for ensemble
ensemble_prediction =
np.mean(list(predictions.values()), axis=0)

# Evaluate the ensemble performance
ensemble_mse = mean_squared_error(y_imputed,
ensemble_prediction)
ensemble_r2 = r2_score(y_imputed,
ensemble_prediction)

print(f"Averaging Ensemble - MSE: {ensemble_mse}, R-
squared: {ensemble_r2}")
```

```
142
143 [8 rows x 506 columns]
144 Best hyperparameters for Linear Regression: {}
145 Linear Regression - Cross-Validation RMSE: 4807493987094.403, Cross-Validation
    R2: -4.358248093168656e+25
146 -----
147 Linear Regression - MSE: 0.3567972677464795, MAE: 0.3645591591735241, R-
    squared: 0.8313722680859739, Adjusted R-squared: 0.6185901425286242
148 -----
149 Averaging Ensemble - MSE: 0.09693058720760624, R-squared: 0.9541891529126637
150 Best hyperparameters for Ridge Regression: {'alpha': 10.0}
151 Ridge Regression - Cross-Validation RMSE: 0.9404219375574672, Cross-Validation
    R2: -0.8914074299466396
152 -----
153 Ridge Regression - MSE: 0.7163920723221295, MAE: 0.6206280253881438, R-squared
    : 0.6614223783722861, Adjusted R-squared: 0.23418976853781426
154 -----
155 Averaging Ensemble - MSE: 0.09563523830784779, R-squared: 0.9548013542010402
156 Best hyperparameters for Lasso Regression: {'alpha': 0.1}
157 Lasso Regression - Cross-Validation RMSE: 1.0152975845700734, Cross-Validation
    R2: -1.0358410669035547
158 -----
159 Lasso Regression - MSE: 0.9949114863899965, MAE: 0.7350715887035172, R-squared
    : 0.5297899323478972, Adjusted R-squared: -0.0635424722205915
160 -----
161 Averaging Ensemble - MSE: 0.09792519579310954, R-squared: 0.9537190860005067
162 Best hyperparameters for SVR: {'C': 1, 'kernel': 'rbf'}
163 SVR - Cross-Validation RMSE: 0.8807387823360168, Cross-Validation R2: -0.
    7810958107585633
164 -----
165 SVR - MSE: 0.4589693954377116, MAE: 0.4146557481507459, R-squared: 0.
    7830841904719806, Adjusted R-squared: 0.5093699769528339
166 -----
167 Averaging Ensemble - MSE: 0.09656035065169855, R-squared: 0.9543641322534213
168 Best hyperparameters for Random Forest Regressor: {'max_depth': None, '

```

File - Project

```
168 n_estimators': 50}
169 Random Forest Regressor - Cross-Validation RMSE: 0.9132126876153362, Cross-
    Validation R2: -0.842168931150432
170 -----
171 Random Forest Regressor - MSE: 0.14607493312023298, MAE: 0.2692731574634776, R
    -squared: 0.9309627990787747, Adjusted R-squared: 0.8438485255971288
172 -----
173 Averaging Ensemble - MSE: 0.09693223440883281, R-squared: 0.9541883744206966
174 Best hyperparameters for Decision Tree Regressor: {'max_depth': 5}
175 Decision Tree Regressor - Cross-Validation RMSE: 1.020996846352605, Cross-
    Validation R2: -1.0472205301796547
176 -----
177 Decision Tree Regressor - MSE: 0.700656151069886, MAE: 0.6189834597437843, R-
    squared: 0.6688594104076014, Adjusted R-squared: 0.2510111851363952
178 -----
179 Averaging Ensemble - MSE: 0.09688216909062473, R-squared: 0.9542120360398292
180 Best hyperparameters for XGBoost Regressor: {'learning_rate': 0.1, 'max_depth
    ': 3, 'n_estimators': 100}
181 XGBoost Regressor - Cross-Validation RMSE: 0.9097185159823049, Cross-
    Validation R2: -0.8360326411742485
182 -----
183 XGBoost Regressor - MSE: 0.5261224563726364, MAE: 0.5462509838439292, R-
    squared: 0.751346648231083, Adjusted R-squared: 0.43758456345534236
184 -----
185 Averaging Ensemble - MSE: 0.09830180553057682, R-squared: 0.9535410945987033
186 Best hyperparameters for KNN Regressor: {'n_neighbors': 7}
187 KNN Regressor - Cross-Validation RMSE: 0.9313758271882367, Cross-Validation R2
    : -0.8753681365879485
188 -----
189 KNN Regressor - MSE: 0.6032714876748544, MAE: 0.5572387754457035, R-squared: 0
    .7148848607010814, Adjusted R-squared: 0.35511363754583747
190 -----
191 Averaging Ensemble - MSE: 0.09694261601587284, R-squared: 0.9541834679177409
192 Best hyperparameters for Gaussian Process Regressor: {}
193 Gaussian Process Regressor - Cross-Validation RMSE: 16.900493907775875, Cross-
    Validation R2: -778.2621035921362
194 -----
195 Gaussian Process Regressor - MSE: 0.03360261696026889, MAE: 0.
    06365532605722445, R-squared: 0.9841189000123955, Adjusted R-squared: 0.
    964079407259957
196 -----
197 Averaging Ensemble - MSE: 0.09752392042011414, R-squared: 0.9539087347510382
198 Best hyperparameters for Bayesian Ridge Regressor: {'alpha_1': 1e-06, 'alpha_2
    ': 0.0001, 'lambda_1': 0.0001, 'lambda_2': 1e-06}
199 Bayesian Ridge Regressor - Cross-Validation RMSE: 0.9409043756658295, Cross-
    Validation R2: -0.8921524045933928
200 -----
201 Bayesian Ridge Regressor - MSE: 0.7049439524033181, MAE: 0.6153278246715537, R
    -squared: 0.6668329312859386, Adjusted R-squared: 0.2464276026841553
202 -----
203 Averaging Ensemble - MSE: 0.09700421645860309, R-squared: 0.954154354626016
204
205 Process finished with exit code 0
206
```

- ENSEMBLE LEARNING IS A MACHINE LEARNING TECHNIQUE THAT COMBINES MULTIPLE MODELS TO CREATE A MORE ACCURATE AND ROBUST MODEL.
- TO IMPROVE ACCURACY: ENSEMBLE MODELS CAN OFTEN OUTPERFORM SINGLE MODELS, ESPECIALLY WHEN THE SINGLE MODELS ARE DIVERSE. THIS IS BECAUSE THE ENSEMBLE MODEL CAN LEARN FROM THE STRENGTHS OF EACH INDIVIDUAL MODEL AND COMPENSATE FOR THEIR WEAKNESSES. TO REDUCE VARIANCE: ENSEMBLE MODELS CAN ALSO HELP TO REDUCE THE VARIANCE OF A MODEL, WHICH MEANS THAT THE MODEL WILL BE LESS LIKELY TO OVERFIT THE TRAINING DATA. THIS IS BECAUSE THE ENSEMBLE MODEL IS AVERAGING THE PREDICTIONS OF MULTIPLE MODELS, WHICH HELPS TO SMOOTH OUT THE NOISE IN THE DATA. TO INCREASE ROBUSTNESS: ENSEMBLE MODELS CAN ALSO BE MORE ROBUST TO CHANGES IN THE DATA THAN SINGLE MODELS. THIS IS BECAUSE THE ENSEMBLE MODEL IS NOT AS SENSITIVE TO THE ERRORS OF ANY INDIVIDUAL MODEL.
- **Averaging Ensemble Observation:**



<u>Regression models</u>	<u>R-squared</u>
➤ <u>LINEAR REGRESSION</u>	➤ <u>0.9544910020588832</u>
➤ <u>RIDGE REGRESSION</u>	➤ <u>0.9544241790991296</u>
➤ <u>LASSO REGRESSION</u>	➤ <u>0.9542165679486984</u>
➤ <u>KNN REGRESSOR</u>	➤ <u>0.9540503777797479</u>
➤ <u>SVR</u>	➤ <u>0.9541140448036425</u>
➤ <u>RANDOM FOREST REGRESSOR</u>	➤ <u>0.9540431206483831</u>
➤ <u>GAUSSIAN PROCESS REGRESSOR</u>	➤ <u>0.9546869179063504</u>
➤ <u>BAYESIAN RIDGE REGRESSOR</u>	➤ <u>0.9546315000197529</u>
➤ <u>DECISION TREE REGRESSOR</u>	➤ <u>0.9542307744311099</u>
➤ <u>XGBOOST REGRESSOR</u>	➤ <u>0.9544810929191012</u>

WE DO OBSERVE THAT THE ENSEMBLE MODELS OUTPERFORMED SINGLE MODELS.

YET, THE IMPORTANT OBSERVATION HERE IS THAT GAUSSIAN PROCESS REGRESSOR (ON ITS SINGLE MODEL) HAS AN R-SQUARED SCORE OF 0.984 WHILE THE ENSEMBLE METHODS R-SQUARED LIES AROUND 0.95.

SO, WE CONCLUDE THAT:

GAUSSIAN PROCESS REGRESSOR IS A BETTER FIT FOR THE PREDICTION MODEL WITH AN R-SQUARED SCORE OF 0.984.

THANK YOU